

# Time-focused density-based clustering of trajectories of moving objects

Margherita D’Auria<sup>1</sup>, Mirco Nanni<sup>2</sup>, and Dino Pedreschi<sup>1</sup>

<sup>1</sup>Dipartimento di Informatica, Università di Pisa  
Via F. Buonarroti 2, 56127 Pisa, Italy  
email: { dauria@cli. | pedre@ } di.unipi.it

<sup>2</sup>ISTI - Institute of CNR  
Via Moruzzi 1 – Loc. S. Cataldo, 56124 Pisa  
email: mirco.nanni@isti.cnr.it

**Abstract.** Spatio-temporal, geo-referenced datasets are growing rapidly, and will be more in the near future, due to both technological and social/commercial reasons. From the data mining viewpoint, spatio-temporal trajectory data introduce new dimensions and, correspondingly, novel issues in performing the analysis tasks. In this paper, we consider the clustering problem applied to the trajectory data domain. In particular, we propose an adaptation of a density-based clustering algorithm to trajectory data based on a simple notion of distance between trajectories. Then, a set of experiments on synthesized data is performed in order to test the algorithm and to compare it with other standard clustering approaches. Finally, a new approach to the trajectory clustering problem, called *temporal focussing*, is sketched, having the aim of exploiting the intrinsic semantics of the temporal dimension to improve the quality of trajectory clustering.

**Note:** The authors are members of the Pisa KDD Laboratory, a joint research initiative of ISTI-CNR and the University of Pisa: <http://www-kdd.isti.cnr.it>.

## 1 Introduction

Spatio-temporal, geo-referenced datasets are growing rapidly, and will be more in the near future. This phenomenon is due to the daily collection of transaction data through database systems, network traffic controllers, web servers, sensor networks. In prospect, an important source is telecommunication data from mobile phones and other location-aware devices – data that arise from the necessity of tracking such wireless, portable devices in order to support their interaction with the network infrastructure. But other than ordinary communication operations, the large availability of these forms of geo-referenced information is expected to enable novel classes of applications, where the discovery of knowledge is the key step. As a distinguishing example, the presence of a large number of location-aware, wireless mobile devices presents a growing possibility to access their tracking logs and reconstruct space-time trajectories of these personal devices and their human companions: trajectories are indeed the traces of moving objects and individuals. These mobile trajectories contain detailed information about personal and vehicular mobile behaviour, and therefore offer interesting practical opportunities to find behavioural patterns, to be used for instance in traffic and sustainable mobility management.

However, spatio-temporal data mining is still in its infancy, and even the most basic questions in this field are still largely unanswered: what kinds of patterns can be extracted from

trajectories? Which methods and algorithms should be applied to extract them? One basic data mining method that could be applied to trajectories is *clustering*, i.e., the discovery of groups of *similar* trajectories.

Spatio-temporal trajectory data introduce new dimensions and, correspondingly, novel issues in performing the clustering task. Clustering moving object trajectories, for example, requires finding out both a proper spatial granularity level and significant temporal sub-domains. Moreover, it is not obvious to identify the most promising approach to the clustering task among the many in the literature of data mining and statistics research; neither it is obvious to choose among the various options to represent a trajectory of a moving objects and to formalize the notion of (dis)similarity (or distance) among trajectories. A brief account of the research in this area is reported in Section 2.

In this context, we precisely address the problem of trajectory clustering. Our basic assumption on source data is that a collection of individual trajectories of moving objects can be reconstructed in an approximated way on the basis of tracking log data left by the objects as they move within the network infrastructure. As an example, a mobile phone that moves among the various cells of the wireless network leaves, during its interactions with the network, a set of triples  $(id, loc, t)$ , each specifying the localization at space  $loc$  and at time  $t$  of the phone  $id$ . Starting from the set of triples for a given object  $id$  is therefore possible, in principle, to approximate a function  $f_{id} : time \rightarrow space$ , which assigns a location to object  $id$  for each moment in a given time interval. We call such a function a *trajectory*, and we concentrate on the problem of clustering a given set of trajectories. This basic assumption is consistent with the form of tracking log data that are (or can be) collected in the wireless network infrastructure; for the sake of concreteness, we could count in our research on the availability of a synthesizer of trajectory data, which has been developed at our laboratory [8] and has been used to create the source trajectory datasets employed in the empirical evaluation of the achieved results.

We address two distinct questions: (i) what is the most adequate clustering method for trajectories, and (ii) how can we exploit the intrinsic semantics of the temporal dimension to improve the quality of trajectory clustering.

Concerning the first problem, we advocate that *density-based clustering*, originally proposed in [2], is particularly well-suited to the purpose of trajectory clustering, given its distinctive features:

- the ability to construct non-spherical clusters of arbitrary shape, unlike the classical  $k$ -means and hierarchical methods,
- the robustness with respect to noise in the data,
- the ability of discovering an arbitrary number of clusters to better fit the source data – like hierarchical methods, but with a considerably lower complexity ( $O(n \log n)$  vs.  $O(n^2)$ ).

All the above are key issues for trajectories: it is likely that trajectories of cars in the urban traffic tend to agglomerate in snake-like, non-convex clusters, that many outlier trajectories should not be included in meaningful clusters but rather considered as noise, and that the number of clusters is unpredictable. Density-based clustering algorithms deal with the above problems by agglomerating objects within clusters on the basis of *density*, i.e., the amount of population within a given region in the space. In our approach, we generalize the spatial notion of distance between objects to a spatio-temporal notion of distance between trajectories, and we thus obtain a natural extension of the density-based clustering technique to trajectories. To analyze the consequences of our approach, we consider a particular density-based clustering

algorithm, OPTICS [2], and propose an empirical comparison with several traditional  $k$ -means and hierarchical algorithms; we show how, on a particular experiment, our density-based approach succeeds in finding the natural clusters that are present in the source data, while all the other methods fail. To some extent, this sort of empirical evidence points out that density-based trajectory clustering yields a better quality output, with respect to the other traditional methods.

The second contribution of this paper is *temporal focusing*. Here, we stress that the temporal dimension plays a key role in trajectory clustering: two trajectories, which are very different considering the whole time interval of their duration, may become very similar if restricted considering a smaller sub-interval – an obvious observation with reference to, e.g., the vehicle trajectories in the urban traffic. It is therefore interesting to generalize trajectory clustering with a focus on the temporal dimension – basically enlarging the search space of interesting clusters by considering the restrictions of the source trajectories onto sub-intervals of time. Our proposed algorithm for temporal focussing is therefore aimed at searching the most meaningful time intervals, which allow to isolate the (density-based) clusters of higher quality.

The plan of the paper follows. In the next Section we briefly discuss related work, while in Section 3 we define trajectories and their distances. In Section 4 we revise density-based clustering and the OPTICS algorithm, while in Section 5 we propose the extension of OPTICS to trajectories and empirically evaluate our proposal. Finally, we propose in Section 6 the temporal focusing methods, together with some preliminary empirical experiment.

## 2 Related Work

In recent years, the problem of clustering spatio-temporal data received the attention of several researchers. Most of the actual work is focused on two kinds of spatio-temporal data: moving objects trajectories (the topic of this paper), such as traffic data, and geographically referenced events, such as epidemiological and geophysical data collected along several years.

**Trajectory clustering.** In one of the first works related to the topic, Ketterlin [14] considers generic sequences (thus modelling trajectories as sequences of points) together with a conceptual hierarchy over the sequence elements, used to compute both the cluster representatives and the distance between two sequences. Nanni, one of the authors of the present paper adapted two classical distance-based clustering methods ( $k$ -means and hierarchical agglomerative clustering) to trajectories [18]. In the first part of the present work, we perform a step in the same direction, by adapting instead the density-based approach to trajectories. An alternative strategy is to apply to trajectories some multidimensional scaling technique for non-vectorial data, e.g., Fastmap [6], which maps a given data space to an Euclidean space preserving (approximately) the distances between objects, so that any standard clustering algorithm for vectorial data can be applied. Other distances can be inherited from the time-series domain, with the implicit assumption that the temporal component of data can be safely ignored and replaced by an order in the collected data values: the most common distance is the Euclidean metric, where each value of the series becomes a coordinate of a fixed-size vector; other approaches, which try to solve problems such as noise and time warping, include the computation of the longest common subsequence of two series (e.g., see [20]), the count of common subsequences of length two [1], the domain-dependent extraction of a single representative value for the whole series [13], and so on. The main drawback of this transformational approach is ad-hoc nature, bound to specific applications. A thoroughly different method, proposed by Gaffney and Smyth [7], is model-based clustering for continuous trajectories, which groups

together objects which are likely to be generated from a common core trajectory by adding Gaussian noise. In a successive work [4] spatial and (discrete) temporal shifting of trajectories within clusters is also considered.

**Spatio-temporal density.** The problem of finding densely populated regions in space-time is conceptually closely related to clustering, and it has been undertaken along several different directions. In [10], a system is proposed to support density queries over a database of uniform speed rectilinear trajectories, able to efficiently discover the spatial locations where moving objects are – or will be – dense. In [9] the computational complexity and approximation strategies for a few motion patterns are studied. In particular, the authors study *flock patterns*, that are defined as groups of at least  $n$  moving objects ( $n$  being a parameter) such that there exists a time interval of width larger than a given threshold where all such objects always lay inside a circle of given radius. A much similar objective is pursued in [11], with an emphasis on efficiency issues. Finally, in [17] an extension of *micro-clustering* to moving objects is proposed, which groups together rectilinear segments of trajectories that lay within a rectangle of given size in some time interval. The objectives of such work, however, are a bit different, being focused on the efficient computation of static clusters at variable time instants. The second contribution of this paper works in a direction similar to micro-clusters discovery. However, in addition to the more general concept of density adopted, in this paper we focus on the global clustering structure of the whole dataset, and not on small groups of objects.

**Event clustering.** A different view of the spatio-temporal clustering problem consists in considering spatially and temporally referenced events, instead of moving objects. In this case, the most basic approach consists in the application of spatial clustering algorithms where time becomes an additional spatial dimension. In addition to that, in literature different approaches have been proposed, mostly driven by specific application domains. Among them, we mention Kuldorff’s spatial scan [16], a well known statistical method developed for epidemiological data which searches spatio-temporal cylinders (i.e., spatial circular shapes which remain still for some time interval) where the rate of disease cases is higher than outside the cylinder, and extensions for considering more flexible square pyramid shapes [12].

### 3 A data model and distance for trajectories

In this paper we consider databases composed by a finite set of spatio-temporal objects. From an abstract point of view, a spatio-temporal object  $o$  is represented by a *trajectory*  $\tau_o$ , i.e., a continuous function of time which, given a time instant  $t$ , returns the position at time  $t$  of the object in a  $d$ -dimensional space (typically  $d \in \{2, 3\}$ ). Formally:  $\tau_o : \mathbb{R}^+ \rightarrow \mathbb{R}^d$ .

In a real-world application, however, trajectories of objects are given by means of a finite set of observations, i.e. a finite subset of points taken from the real continuous trajectory. Moreover, it is reasonable to expect that observations are taken at irregular rates within each object, and that there is not any temporal alignment between the observations of different objects. This calls for an approximate reconstruction of the original trajectory. In this paper, we employ the model used in [19], where the objects are assumed to move in a piecewise linear manner. Namely, an object moves along a straight line with some constant speed till it changes the direction and/or speed. Such model essentially corresponds to the well known parametric 2-spaghetti approach [3].

In this work, we are interested in distances that describe the similarity of trajectories of objects along time and therefore are computed by analyzing the way the distance between the objects varies. More precisely, we restrict to consider only pairs of *contemporary* instantiations

of objects, i.e., for each time instant we compare the positions of the objects at that moment, thus aggregating the set of distance values obtained this way. This implies, in particular, that we exclude subsequence matching and other similar operations usually adopted in the time series field, as well as solutions that try to align – in time and/or in space – shifted trajectories.

The distance between trajectories adopted in this paper is computed in a most natural way, as the average distance between objects, that is to say:

$$D(\tau_1, \tau_2)|_T = \frac{\int_T d(\tau_1(t), \tau_2(t))dt}{|T|},$$

where  $d()$  is the Euclidean distance over  $\mathcal{R}^2$ ,  $T$  is the temporal interval over which trajectories  $\tau_1$  and  $\tau_2$  exist, and  $\tau_i(t)$  ( $i \in \{1, 2\}$ ) is the position of object  $\tau_i$  at time  $t$ . We notice that such a definition requires a temporal domain common to all objects, which, in general, is not a hard requirement. From a conceptual viewpoint, moreover, in order to compute  $D()$  we need to compute the infinite set of distances for each  $t \in T$  (e.g., in and, afterward, to aggregate them. However, due to the (piece-wise) linearity of our trajectories, it can be shown [18] that  $D()$  can be computed as a finite sum by means of  $O(n_1 + n_2)$  Euclidean distances,  $n_1$  and  $n_2$  being the number of observations respectively available for  $\tau_1$  and  $\tau_2$ . Moreover, such distance is a metric, thus allowing the use of several indexing techniques that help to improve performances in several applications, including clustering.

## 4 Density-based clustering and OPTICS

In this section we briefly review the principles of density-based clustering, summarizing the motivations which led us to adopt this approach for trajectories, and describe the OPTICS algorithm.

### 4.1 Density-based Clustering

The key idea of density-based clustering algorithms, and OPTICS in particular, is that for each object in some cluster the neighborhood of a given radius  $\epsilon$  has to contain at least a minimum number  $n_{pts}$  of objects, i.e., the cardinality of the neighborhood has to exceed a given threshold. For that reason, such algorithms are naturally robust to problems such as noise and outliers, since they usually do not significantly affect the overall density distribution of data. This is an important feature for several real world applications, such as all those that work with data sources having some underlying random (unpredictable) component – e.g., data obtained by observing human behaviour – or that collect data by means of not-completely reliable methods – e.g., low-resolution sensors, sampled measures, etc. Trajectory data usually suffer of both the mentioned problems, so noise tolerance is a major requisite.

Moreover, in typical applications dealing with moving objects, such as traffic analysis or the study of PDAs usage and mobility, any strict constraint to the shape of clusters would be a strong limitation, so  $k$ -means and other spherical-shape clustering algorithms could not be generally applied (e.g., large groups of cars moving along the same road would form a "snake"-shaped cluster, not a spherical one). It is worth noting that the output of OPTICS, the *reachability plot*, described in the following, is an intuitive, data-independent visualization of the cluster structure of data, that yields valuable information for a better comprehension of the data and that is (also) used to assign each object to its corresponding cluster or to noise, respectively.

## 4.2 OPTICS

In the following, we will shortly introduce the definitions underlying OPTICS, i.e., *core objects* and the *reachability-distance* of an object  $p$  w.r.t. a predecessor object  $o$ , and briefly describe how the algorithm works by means of a small example.

An object  $p \in D$  is called a *core object* if the neighborhood around it is a dense region, and therefore  $p$  should definitely belong to some cluster and not to the noise. More formally:

**Definition 1 (Core object).** Let  $p \in D$  be an object,  $\varepsilon$  a distance threshold and  $N_\varepsilon(p)$  the  $\varepsilon$ -neighborhood of  $p$ , i.e., the set of points  $\{x \in D | d(p, x) \leq \varepsilon\}$ . Then, given a parameter  $n_{pts} \in \mathcal{N}$ ,  $p$  is a core object if:  $|N_\varepsilon(p)| \geq n_{pts}$ .

Based on core objects, we have the following:

**Definition 2 (Reachability-distance).** Let  $p \in D$  be an object,  $o \in D$  a core object,  $\varepsilon$  a distance threshold and  $N_\varepsilon(o)$  the  $\varepsilon$ -neighborhood of  $o$ . Denoting with  $n\text{-distance}(p)$  the distance from  $p$  to its  $n$ -th neighbor in order of proximity ( $n \in \mathcal{N}$ ), and given a parameter  $n_{pts} \in \mathcal{N}$ , the reachability-distance of  $p$  with respect to  $o$  is defined as

$$\text{reach-}d_{\varepsilon, n_{pts}}(p, o) = \max\{n_{pts}\text{-distance}(o), d(o, p)\}$$

The reachability-distance of  $p$  w.r.t.  $o$  is essentially their distance, excepted when  $p$  is *too close*, in which case such distance is *normalized* to a suitable value. OPTICS works as follows: initially a random object  $p_0$  is chosen; then, at each iteration  $i$ , the next object  $p_i$  chosen from  $D$  is that with the smallest reachability-distance w.r.t. all the already visited core objects; the process is repeated until all objects in  $D$  have been considered.

The whole process is summarized by means of a *reachability plot*: on the horizontal axis are represented the objects in their visit ordering  $0, \dots, |D| - 1$ , and on the vertical axis, for each  $i$  the reachability-distance corresponding to  $p_i$  is plotted. The sequence  $\langle p_0, \dots, p_{|D|-1} \rangle$  is also called a *cluster-reordering* of the objects in  $D$ .

Intuitively, the reachability-distance of a point  $p_i$  corresponds to the minimum distance from the set of its predecessors  $p_j, 0 \leq j < i$ . As a consequence, a high value of reachability-distance approximatively means a high distance from all other objects, and therefore indicates a *rarefied* area. Then, clusters, i.e., dense areas, are represented as valleys in the reachability plot. Figure 1 shows a sample execution of OPTICS on a toy dataset with two clusters, and the corresponding values on the reachability plot. We observe that the jump from point 9 to point 10, belonging to different clusters, corresponds to a peak in the reachability plot.

From the reachability plot we can easily obtain a partitioning of the data into a set of clusters, plus noise: we simply need to choose a threshold to separate clusters, expressed in terms of the maximum value of the reachability distance allowed within clusters. Such value, denoted  $\varepsilon'$  and set by the user, is used to separate the objects into peaks and valleys: the former will be considered noise, the latter as clustered objects. In the example in Figure 1, setting  $\varepsilon' = 0.1$ , we would obtain two clusters: objects 1-9 and 11-18.

## 5 Extending OPTICS to trajectories

In order to define similarity measures for OPTICS over complex domains, its inventors proposed and tested a few solutions, classified into two classes: feature-based models and direct

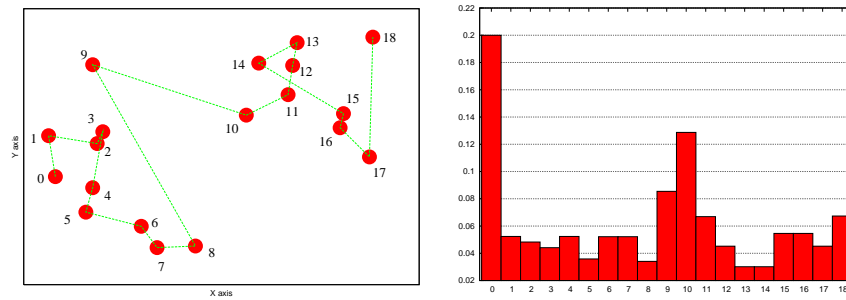


Fig. 1. Sample run of OPTICS with the resulting reachability plot

geometric models [15]. In the first case, the basic idea is to use a feature extraction function that maps the objects onto vectors in an appropriate multidimensional feature space. The similarity of two objects is then defined as their proximity in the feature space, and the closer their feature vectors are located, the more similar the objects are considered. In the second class of models, the distance is defined by directly using the geometry of objects. Examples include the computation of the volume of the geometric intersection of the compared objects. In this paper, we follow an approach similar to the latter, since trajectories are compared and clustered by means of the spatio-temporal distance described in Section 3. An adequate efficiency can moreover be preserved by adopting some indexing tool for generic metrics, such as M-tree [5]. However, it is less obvious to verify that OPTICS deliver high quality trajectory clusters. In the next section, we provide some evidence that OPTICS finds dense and well defined trajectory clusters, with a good tolerance to noise and outliers.

### 5.1 OPTICS vs. Hierarchical and K-means algorithms

We considered a synthetic test dataset randomly generated with the CENTRE trajectory generator [8], composed of 250 trajectories organized into four natural clusters plus noise. The dataset is depicted in Figure 2(a), where horizontal coordinates represent spatial positions and the vertical one represents time. Each cluster contains objects that move towards some defined direction. The objective of these experiments is to evaluate the behaviour of OPTICS on the trajectory data domain, as compared to other classical, general purpose clustering algorithms. We applied to the dataset the  $k$ -means algorithm and a standard hierarchical agglomerative algorithm in three versions: single-linkage, complete-linkage and average-linkage. The same dataset was processed by the trajectory version of OPTICS. All algorithms were configured to find four clusters (i.e., the expected number of natural clusters in the data). For the output of each algorithm, the average *purity* of the clusters (i.e., percentage of objects in the cluster that belonged to the corresponding real cluster) and the average of their *coverage* (i.e., percentage of objects of the real cluster that appear in the cluster found). The results are discussed below:

- K-means yields a 100% coverage but only a 53.7% purity. That is due to the fact that it merges together two clearly distinguished clusters, since it is sensible to noise and outliers. In fact, a whole resulting cluster is composed only of noisy objects.
- The hierarchical single-linkage algorithm yields a 100% coverage but only a 52.0% purity. Indeed, it is very sensitive to the chaining effect – i.e., the fact of collapsing far clusters

due to a thin, yet continuous, line of objects linking them – and therefore it merges two dense and well separated clusters because of noise and the closeness of their borders.

- The complete-linkage algorithm yields a 99% coverage but only a 50.0% purity: it tends to form clusters with equal diameter, and so, due to the presence of noise that biases such size, two of the natural clusters are again merged together.
- The average-linkage algorithm yields a 100% coverage but only a 50.5% purity: it keeps clusters with balanced average intra-cluster distance, which usually results in a behavior similar to the complete-linkage case. Again, a pair of natural clusters is merged.
- Finally, trajectory-OPTICS yields a 93.5% coverage and also a 99.0% purity, i.e., the best trade-off between the two measures. As we can see from the resulting reachability plot (Figure 2(b)), trajectory-OPTICS correctly finds the four natural clusters, which can be easily isolated by selecting a proper value for the  $\epsilon$  parameter ( $\epsilon = 24$  in our example, but any value that crosses the three central protrusions yield the same result).

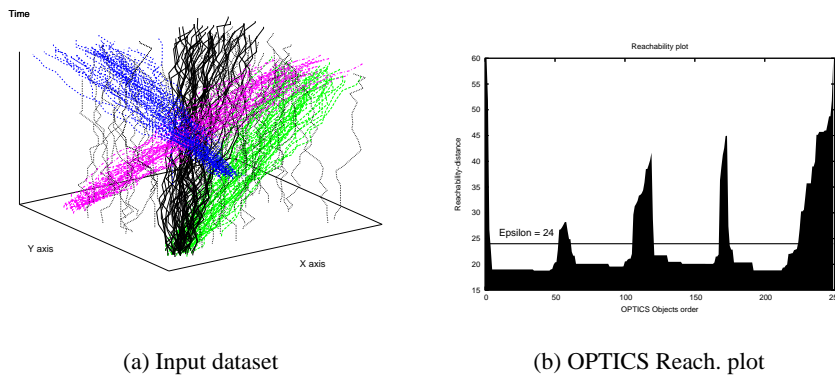


Fig. 2. A synthesized dataset (a) and the corresponding reachability plot for OPTICS

## 6 Temporal Focusing

The approach to trajectory clustering presented above treats trajectories as unique, indivisible elements, and tries to group together those moving objects that globally move in a similar way, "smoothing" the effect of any sporadic divergence in their movement. However, such global trajectory clustering may sometime be misleading and yield counter-intuitive results. In particular, it might keep separated objects that moved together for a significant amount of time, while gathering objects that constantly keep a significant distance between them.

From real world experience, we learnt that not all time intervals have the same importance. A meaningful example is urban traffic: in rush hours a large quantity of people move from home to work and viceversa, or, more generally, from/to largely shared targets. Therefore, we can expect that the sheer size of the population sample will make it possible for groups of individuals having similar destinations to clearly emerge from traffic data to form compact clusters. In quiet periods of the day, on the contrary, we expect to mainly observe static individuals, whose distribution on the territory is more driven by the geographical population density than by collective motion behaviors. This is a general problem not limited to urban traffic, and, while in this sample context some interesting hours of the day for cluster analysis can be guessed – e.g., typical morning rush hours –, in other, less understood cases and domains

it might be not possible to fix a priori criteria for choosing the right period of time. In these cases, some automatic mechanism to discover the most interesting intervals of time should be applied. In what follows we formalize the problem mentioned above, and suggest a solution. We anticipate that we will consider only single intervals of time in our search problem, thus not taking into account more complex patterns of time, such as periodical patterns or irregular sets of disjoint time intervals.

## 6.1 Problem setting

As discussed above, there may exist time segments where the clustering structure of our moving objects dataset is clearer than just considering the whole trajectories. In order to discover such clustering structure, then, we should provide a method for locating the right time interval, and focus the clustering process on the segments of trajectories that lay in that interval, ignoring the remaining parts.

Our general approach consists in the following: the trajectory-OPTICS algorithm is modified to compute distances between trajectories focusing on any time interval specified by the user; therefore, we (hypothetically) apply such algorithm to each possible time interval, evaluate each result obtained and determine the best one. This entails solving the following optimization problem:

$$\arg \max_{\theta} Q(D, \theta)$$

where  $D$  is the input dataset and  $Q$  is a quality function that measures the goodness of the clustering results obtained with the parameters  $\theta$ . The set of parameters  $\theta$  can contain, in a general setting, several different components, e.g.: basic parameters for OPTICS, a temporal window, a spatial granularity, etc. In this work we will focus on time windows, so  $\theta = \langle \varepsilon, \varepsilon', n_{pts}, I \rangle$ ,  $\varepsilon$ ,  $\varepsilon'$  and  $n_{pts}$  being the already mentioned general parameters for OPTICS and  $I$  being the time window we are focusing on. Moreover, since OPTICS is not very sensible to its input parameters  $n_{pts}$  and  $\varepsilon$ , we can assume that they are set before approaching the optimization problem, so that the only variables of the problem remain  $I$  and  $\varepsilon'$ .

We observe that the problem introduced above can be seen as a subspace clustering problem. In particular, if we assume that time is discrete, and therefore the time interval  $I$  can be reduced to a finite sequence of  $N$  time points, trajectories can be seen as  $2N$ -dimensional objects, and our objective is to discover  $2N'$  contiguous dimensions ( $N' \leq N$ ) that optimize our quality function. However, in all works of the subspace clustering literature, the dimensions are not related to each other by any specific semantics, differently from the temporal semantics that underlies trajectories. In particular, distances that are additive w.r.t. dimensions<sup>1</sup> are usually applied, making density of regions a monotonic function w.r.t. the dimensions (i.e., dense regions on some  $N$ -dimensional space always correspond to dense regions on any of its  $N - 1$ -dimensional subspaces). On the contrary, when dealing with trajectories we have to deal with a semantics of time, that significantly modifies the problem. The most direct and relevant effects are the following: (i) a notion of contiguity is defined, that has to be preserved when selecting subspaces; and (ii) a distance between objects is given, that is not additive w.r.t. dimensions. As a consequence, actual subspace clustering techniques are not applicable in our

<sup>1</sup> I.e., distances between  $N$ -dimensional objects can be written as a sum of  $N$  contributions independently computed on each dimension. E.g., it holds for Euclidean distances, since  $d^2(a, b) = \sum_{i=1}^N (a_i - b_i)^2$ .

case. In what follows, we propose an heuristic solution to this new variant of the subspace clustering problem.

## 6.2 Quality measure

The first issue to solve, now, is the definition of a quality function  $Q$ , which can provide a good criterion for deciding whether a clustering result is better than another. Any suitable definition should take into account the nature of clusters we obtain. In particular, since we are working with density-based tools, typical dispersion measures cannot be applied, because we can have good non-spherical clusters, which, in general, are not compact.

In the density-based setting, the standard *high intra-cluster vs. low inter-cluster similarity* principle could be naturally translated into a *high-density clusters vs. low-density noise* rule. Highly dense clusters can be considered interesting per se, while having a rarefied noise means that clusters are clearly separated. Put together, these two qualities seem to reasonably qualify a *good* (density-based) clustering result.

The reachability plot returned by the OPTICS algorithm contains a summary of the information on data density we need. Therefore, we can simplify the computation of the  $Q$  measure by deriving it from the corresponding reachability plot, since density at each point can be estimated by the corresponding reachability-distance.

**Definition 3.** *Let  $D$  be an input dataset of trajectories,  $I$  a time interval and  $\epsilon'$  a density threshold parameter. Then, the average reachability,  $R_{D,I,\epsilon'}$ , is defined as the average reachability-distance of non-noise objects:  $R_{D,I,\epsilon'} = \text{avg}\{r | \langle p_0, \dots, p_{|D|-1} \rangle$  is OPTICS cluster reordering  $\wedge r = \text{reach-d}_{\epsilon, n_{pts}}(p_i) \wedge r \leq \epsilon'\}$ . When clear from the context, average reachability will be denoted as  $R_C$  (reachability of clustered objects). When no  $\epsilon'$ -cut is specified (i.e.,  $\epsilon' = \infty$ ), average reachability will also be denoted as  $R_G$  (global reachability).*

**Definition 4.** *Let  $D$  be an input dataset of trajectories,  $I$  a time interval and  $\epsilon'$  a density threshold parameter. Then, we define the  $Q_1$  quality measure as follows:*

$$Q_1(D, I, \epsilon') = -R_{D,I,\epsilon'}.$$

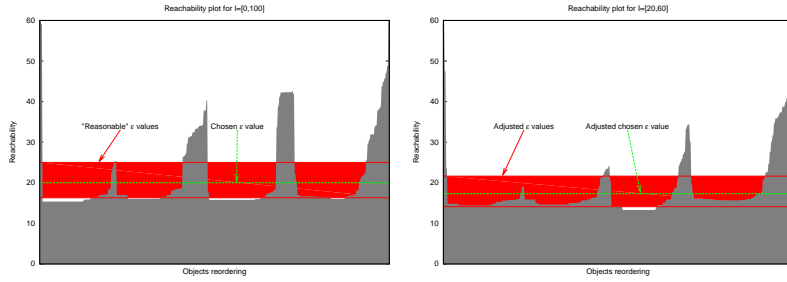
## 6.3 Self-tuning of parameters

The search for an optimal time interval requires to iteratively analyze different views of the input trajectory dataset. However, each view can result in a quite different reachability plot, with peaks and valleys of highly variable size. As a consequence, it is not possible to determine a single value of the  $\epsilon'$  that is valid for all time intervals. Since it is not reasonable to ask the user to choose a suitable value for each interval analyzed, we should design an automated method for the purpose. We describe a solution to the problem that requires the user to specify an initial  $\epsilon'_0$  value for the largest time interval, and then automatically adapts such value to the smaller intervals.

In general, the optimal  $\epsilon'$  value for an interval  $I$  depends on several variables. We consider a simple solution and assume that densities follow some general trend, so that a global rescaling factor for density values can be computed. Such rescaling can be applied to  $\epsilon'_0$  to obtain the actual  $\epsilon'$ . The overall density for a given time interval can be estimated in several ways, the simplest being the average reachability-distance  $R_G$ . Adopting such measure and a linear rescaling, we have the following value for  $\epsilon'$ :

$$\varepsilon'(D, I, \varepsilon'_0) = \frac{R_G}{R_G^0} \varepsilon'_0$$

In order to test the reliability of the  $\varepsilon'$  re-scaling method proposed above, we performed an experiment on the trajectory dataset considered in Section 5.1. We assumed that for each reachability plot it is possible to define an interval for the acceptable values of the  $\varepsilon'$  parameter, and that the specific  $\varepsilon'$  value that a generic user would choose follows a random uniform distribution over the above mentioned interval. Figure 3 (left) depicts an example of reasonable values over the largest time interval (0 – 100), represented by a dark band, with an example of chosen value for  $\varepsilon$ . The same Figure (right) shows also their mapping to a smaller interval (20 – 60).



**Fig. 3.** Example of  $\varepsilon$  adaptation from  $I = [0, 100]$  to  $I = [20, 60]$

In our experiment, We considered 50 samples that include time intervals of all sizes and of variable position, i.e., not biased towards initial or final intervals. Then, for each time interval  $I$  considered, we compute the probability that a value chosen on the 0 – 100 time interval is *adjusted* to a reasonable value on  $I$ , i.e., a value that lays within corresponding interval of reasonable values. Then, we compute the average probability of success, over all cases. The result is the following: given a random interval  $I$  and an  $\varepsilon$  value randomly chosen on the reachability plot of the 0 – 100 time interval, in the 80.5% of the cases such  $\varepsilon$  is mapped on  $I$  to a reasonable value. Such estimate shows that, in spite of its simplicity, the proposed rescaling method is empirically reliable.

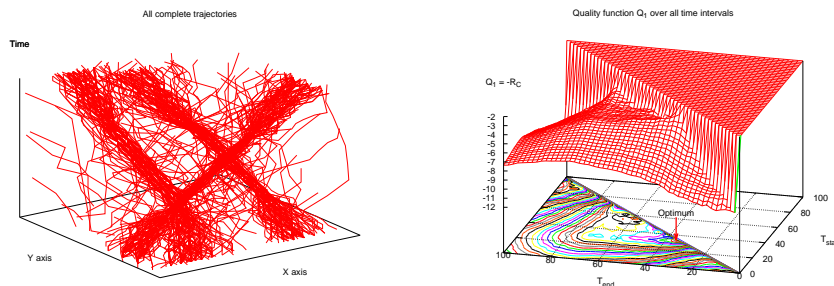
#### 6.4 Evaluating and improving $Q_1$

In order to accept  $Q_1$  as a suitable quality measure, an important property to verify is its *non-shrinking* behavior w.r.t. interval size, i.e., the fact that not all time intervals  $I$  contain at least a subinterval  $I'$  having a higher quality measure. Formally, we would like to refute the following property:  $\forall I. \exists I' \subset I : Q_1(I') > Q_1(I)$ . In fact, such property would imply that the optimal time interval for  $Q_1$  is always a minimal-size interval, thus trivializing the problem and yielding uninteresting results. Fortunately, it is easy to find a counterexample:

*Example 1.* Let consider a set of trajectories over a time interval that is formed by exactly two time units: in the first unit trajectories are very dense, while in the second one they become

more rarefied, although not beyond the given density threshold. Now add a trajectory that is always distant from all the others, but in the first time interval is close enough to them to be part of the same cluster, while in the second one immediately gets extremely far from them. Then, focusing on the first interval, all trajectories belong to the same cluster but, due to the single outlier trajectory described above, it has only a medium density. In the second interval, the outlier trajectory becomes noise, but the remaining ones are a bit rarefied, so density is not high here. On the overall interval, we get all the positive aspects of both the sub-intervals: the outlier trajectory is considered as noise, and so the first sub-interval yields a very high density that is able to lift the overall density beyond the limits of the two sub-intervals. As a consequence, the larger interval has a better quality than its sub-intervals.

The second step in validating the quality measure is a test against a dataset. In Figure 4, on the left, a set of 320 trajectories is shown. Such dataset was generated by using the already mentioned CENTRE system, and contains (in addition to noise) three groups of trajectories that tend to move together within a central time interval, and tend to spread across a wider area in the other intervals. Trajectories are defined over the  $[0, 100]$  time interval, that is discretized into 50 time units of 2 seconds each. In figure 4, on the right, a plot of the  $Q_1$  measure over all time intervals is depicted: the two horizontal axes represent the boundaries of time intervals, so that only the lower-left part of the plane – where the lower bounds of intervals are smaller than upper bounds – is significant. Below the 3D plot, the corresponding (plane) contour plot gives a bi-dimensional view of the same value, with the optimal interval pointed by an arrow. The plot shows that the optimum is located along the diagonal, which represents the intervals with the smallest size. In general, there is a strong bias towards small intervals, even though we can notice the presence of several central regions in the plot – corresponding to larger intervals – with  $Q_1$  values close to the optimum.



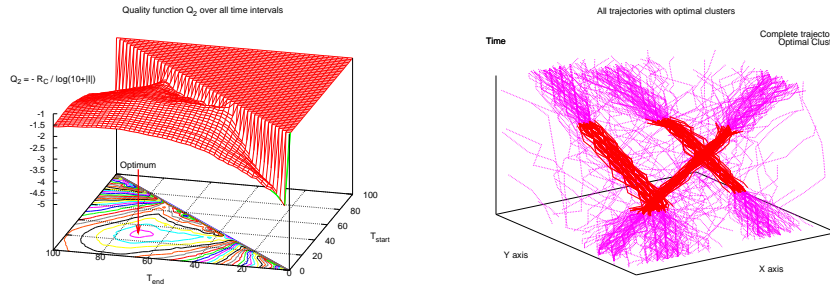
**Fig. 4.** Sample dataset and corresponding  $Q_1$  plot

When only very small variations of  $Q_1$  are observed, it is reasonable to prefer larger intervals, since they are more informative. We can do that by slightly promoting larger intervals directly in the quality measure. A simple solution consists in adding a factor that increases the quality function as the interval size grows, but only very slowly, in order to avoid excessive bias. To this purpose, we adopt the following variation of the  $Q_1$  measure:

$$Q_2(D, I, \epsilon') = Q_1(D, I, \epsilon') / \log_{10}(10 + |I|)$$

In Figure 5 we can see the plot corresponding to the new  $Q_2$  measure. As we can notice, the optimum value (pointed by the arrow) is now in the central region, on the  $[30, 72]$  interval,

which means that the small correction introduced is sufficient to discover significantly large dense intervals. On the righthand plot we see the resulting clusters, where the segments of the clustered trajectories contained in the optimal time interval are emphasized: the densest segments of the natural clusters present in the data are clearly discovered, leaving out the portions of trajectory where the objects start/end dispersing.



**Fig. 5.**  $Q_2$  plot on sample dataset and corresponding output clusters

## 6.5 Searching strategies

The basic method for finding the time interval that maximizes  $Q_2$  is an exhaustive search over all possible intervals. Such approach is obviously very expensive, since it adds to the complexity of OPTICS a quadratic factor w.r.t. the the maximum size of time intervals, expressed in time units. More precisely, given a granularity of time  $\tau$ , the global interval  $I_0$  is divided into  $N_t = |I_0|/\tau$  time units. Then, the cost of an exhaustive search is  $O(N_t^2 \cdot OPTICS(n)) = O(N_t^2 n \log n)$ .

A natural alternative to the exhaustive search of a global optimum is the search of local one, adopting some greedy search paradigm. We considered a standard hill climbing approach, where a randomly chosen solution (i.e., time interval) is iteratively modified, trying to improve the quality function at each step. We follow the procedure described below, where  $\tau$  is the chosen temporal granularity and  $I_0$  the largest interval:

1. Choose an initial random time interval  $I \subseteq I_0$ ;
2. Let  $I' = \arg \max_{T \in Neigh_I} Q_2(D, T, \epsilon')$ , where  $Neigh_I = \{T \in \{[T_s \pm \tau, T_e], [T_s, T_e \pm \tau]\} | T \subseteq I_0\}$  and  $I = [T_s, T_e]$ ;
3. If  $Q_2(D, I', \epsilon') > Q_2(D, I, \epsilon')$  then let  $I := I'$  and return to step 2; otherwise stop.

The total cost of the search procedure is  $O(n_{iter} n \log n)$ , where the number of iterations  $n_{iter}$  has a worst-case value of  $O(N_t^2)$ , but usually assumes much smaller values, linear in  $N_t$ . Executing the greedy search over our sample dataset from all possible starting points, we obtained the exact global optimum in the 70.7% of cases, which provides the success probability when seed points in the search strategy are chosen in a uniformly random way. Therefore, on this sample data we reach a high success probability with a small number of trials (with 5 runs we have over the 99.8% probability of finding the global optimum on this dataset). On the average, each run required around 17 steps and 49 OPTICS invocations – corresponding to less than 4% of the invocations required by the exhaustive search – thus keeping computational costs within reasonable limits.

## 7 Conclusions

In this paper we developed a density-based clustering method for moving objects trajectories, aimed at properly exploiting the intrinsic temporal semantics to the purpose of discovering interesting time intervals, where (when...) the quality of the achieved clustering is optimal. Future research includes, on one hand, a refinement of the general method and a vast empirical evaluation over various real-life datasets, mainly aimed at consolidating (and possibly correcting) the preliminary results shown in this work, and, on the other hand, a deeper integration between the underlying clustering engine and the search method.

## References

1. Rakesh Agrawal, King-Ip Lin, Harpreet S. Sawhney, and Kyuseok Shim. Fast similarity search in the presence of noise, scaling, and translation in time-series databases. In *VLDB*, pages 490–501, 1995.
2. M. Ankerst, M. Breunig, H.-P. Kriegel, and J. Sander. Optics: Ordering points to identify the clustering structure. In *Proc. ACM SIGMOD Int. Conf. on Management of Data (SIGMOD'99)*. ACM Press, 1999.
3. J. Chomicki and P.Z. Revesz. Constraint-Based Interoperability of Spatiotemporal Databases. *GeoInformatica*, 3(3):211–243, 1999.
4. Darya Chudova, Scott Gaffney, Eric Mjolsness, and Padhraic Smyth. Translation-invariant mixture models for curve clustering. In *KDD '03: Procs. of ACM SIGKDD*, pages 79–88. ACM Press, 2003.
5. Paolo Ciaccia, Marco Patella, and Pavel Zezula. M-tree: An efficient access method for similarity search in metric spaces. In *VLDB'97*, pages 426–435. Morgan Kaufmann Publishers, Inc., 1997.
6. C. Faloutsos and K.-I. Lin. Fastmap: a fast algorithm for indexing of traditional and multimedia databases. In *SIGMOD Conf.*, pages 163–174. ACM, 1995.
7. S. Gaffney and P. Smyth. Trajectory clustering with mixture of regression models. In *KDD Conf.*, pages 63–72. ACM, 1999.
8. F. Giannotti, A. Mazzoni, S. Puntoni, and C. Renso. Synthetic generation of cellular network positioning data. Technical report, ISTI-CNR, 2005.
9. Joachim Gudmundsson, Marc J. van Kreveld, and Bettina Speckmann. Efficient detection of motion patterns in spatio-temporal data sets. In *GIS*, pages 250–257, 2004.
10. Marios Hadjieleftheriou, George Kollios, Dimitrios Gunopulos, and Vassilis J. Tsotras. On-line discovery of dense areas in spatio-temporal databases. In *Proceedings of SSTD'03*, 2003.
11. San-Yih Hwang, Ying-Han Liu, Jeng-Kuen Chiu, and Ee-Peng Lim. Mining mobile group patterns: A trajectory-based approach. In *Proceedings of PAKDD'05*, 2005. To appear.
12. Vijay S. Iyengar. On detecting space-time clusters. In *KDD*, pages 587–592, 2004.
13. Konstantinos Kalpakis, Dhiral Gada, and Vasundhara Puttagunta. Distance measures for effective clustering of arima time-series. In *ICDM*, pages 273–280, 2001.
14. A. Ketterlin. Clustering sequences of complex objects. In *KDD Conf.*, pages 215–218. ACM, 1997.
15. Hans-Peter Kriegel, Stefan Brecheisen, Peer Kröger, Martin Pfeifle, and Matthias Schubert. Using sets of feature vectors for similarity search on voxelized cad objects. In *SIGMOD Conference*, pages 587–598, 2003.
16. M. Kulldorff. A spatial scan statistic. *Comm. in Statistics: Theory and Methods*, 26(6):1481–1496, 1997.
17. Yifan Li, Jiawei Han, and Jiong Yang. Clustering moving objects. In *KDD*, pages 617–622, 2004.
18. M. Nanni. *Clustering Methods for Spatio-Temporal Data*. PhD thesis, CS Dept., Univ. of Pisa, 2002.
19. S. Saltenis, C. S. Jensen, S. T. Leutenegger, and M. A. Lopez. Indexing the positions of continuously moving objects. In *Procs. of ACM SIGMOD*, pages 331–342. ACM, 2000.
20. Michail Vlachos, Dimitrios Gunopulos, and George Kollios. Discovering similar multidimensional trajectories. In *ICDE*, pages 673–684, 2002.